



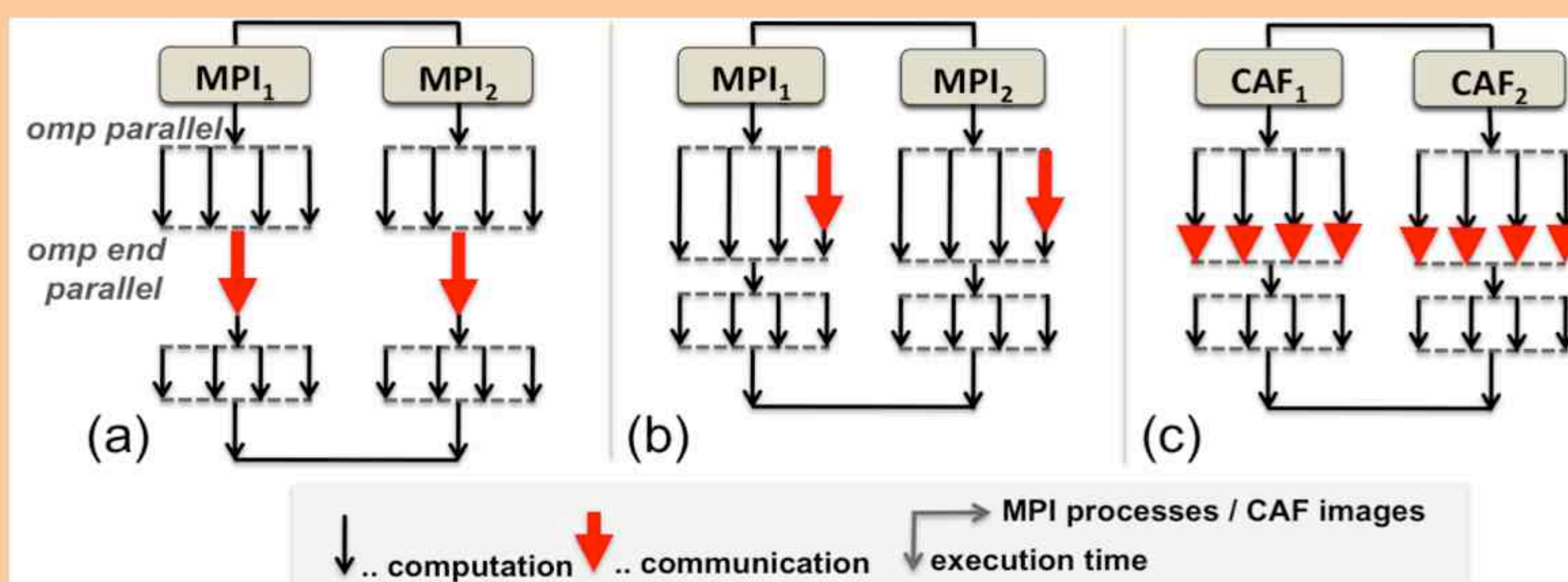
# Global Address Space Communication Techniques for Gyrokinetic Fusion Applications on Ultra-Scale Platforms, *SC'11 best paper finalist*



R Preissl<sup>1</sup>, N Wichmann<sup>2</sup>, B Long<sup>2</sup>, J Shalf<sup>1</sup>, S Ethier<sup>3</sup> and A Koniges<sup>1</sup>

<sup>1</sup>Berkeley Lab (USA), <sup>2</sup>CRAY Inc. (USA), <sup>3</sup>Princeton Plasma Physics Laboratory (USA)

A majority of the work was supported by the Petascale Initiative in Computational Science at NERSC with additional support by the Cray Center of Excellence at NERSC. Additionally, we are grateful for interactions with Nicholas Wright (NERSC) and for the extended computer time as well as the valuable support from NERSC. This work was supported by the Director, Office of Science, Advanced Scientific Computing Research, of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.



## PROGRAMMING MODELS:

- Classical hybrid MPI/OpenMP programming model
- Extension of classical hybrid model where MPI thread teams for work distribution and collective MPI function calls
- The advanced hybrid PGAS/OpenMP algorithm builds on the strategy of communicating threads, but allows ALL OpenMP threads per team to make communication calls to the thread-safe PGAS communication layer

```

1  ! (1) Prepost receive requests
2  do i=1, nr_dests
3    MPLIRECV(recv_buf(i), i, req(i), tor_comm, ...)
4  enddo
5
6  ! (2) compute shifted particles and fill buffer
7  !$omp parallel
8  pack(p_array, shift, holes, send_buf)
9
10 ! (3) Send of particles to destination process
11 do j=1, nr_dests
12   MPLISEND(send_buf(j), j, req(j+i), tor_comm, ...)
13 enddo
14 MPLWAITALL(2*nr_dests, req, ...)

```

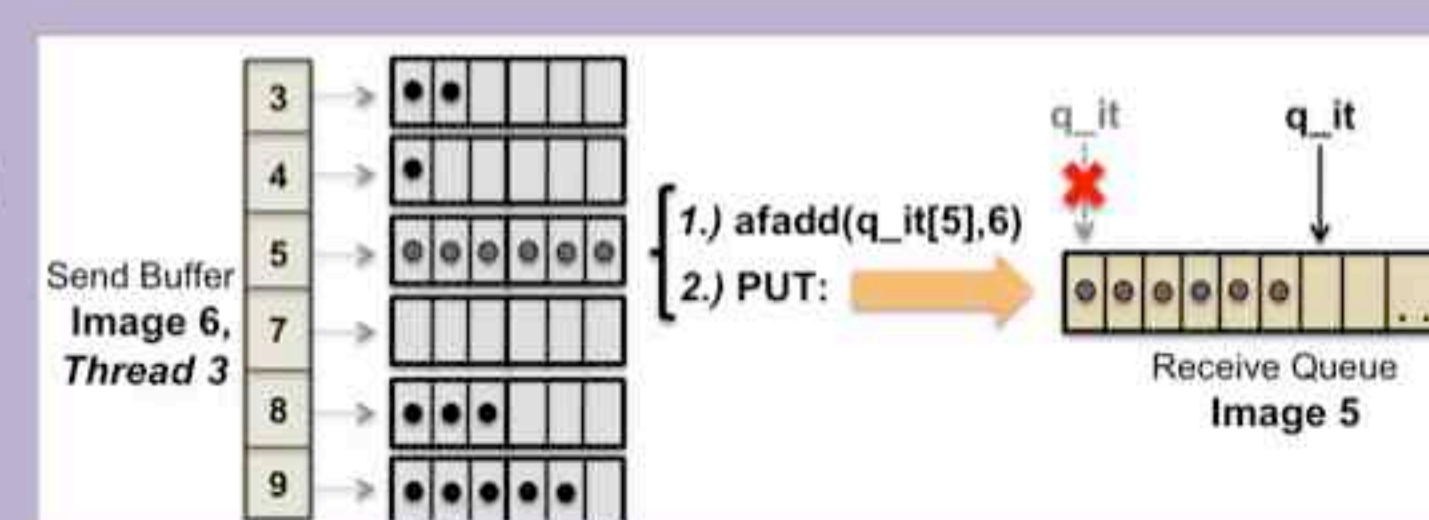
```

1  ! (1) compute shifted particles and fill the
2  ! receiving queues on destination images
3  !$omp parallel do schedule(dynamic, p_size/100)
4  !$omp private(s_buf, buf_cnt) shared(recvQ, q_it)
5  do i=1, p_size
6    dest=compute_destination(p_array(i))
7    if(dest.ne.local_toroidal_domain) {
8      holes(shift++)=i
9      s_buf(dest, buf_cnt(dest))=p_array(i)
10     if(buf_cnt(dest).eq.sb_size) {
11       q_start=afadd(q_it[dest], sb_size)
12       recvQ(q_start:q_start+sb_size-1)[dest] &
13         =s_buf(dest, 1:sb_size)
14       buf_cnt(dest)=0 } }
15   enddo
16
17 ! (2) shift remaining particles
18 empty_s_buffers(s_buf)
19 !$omp end parallel

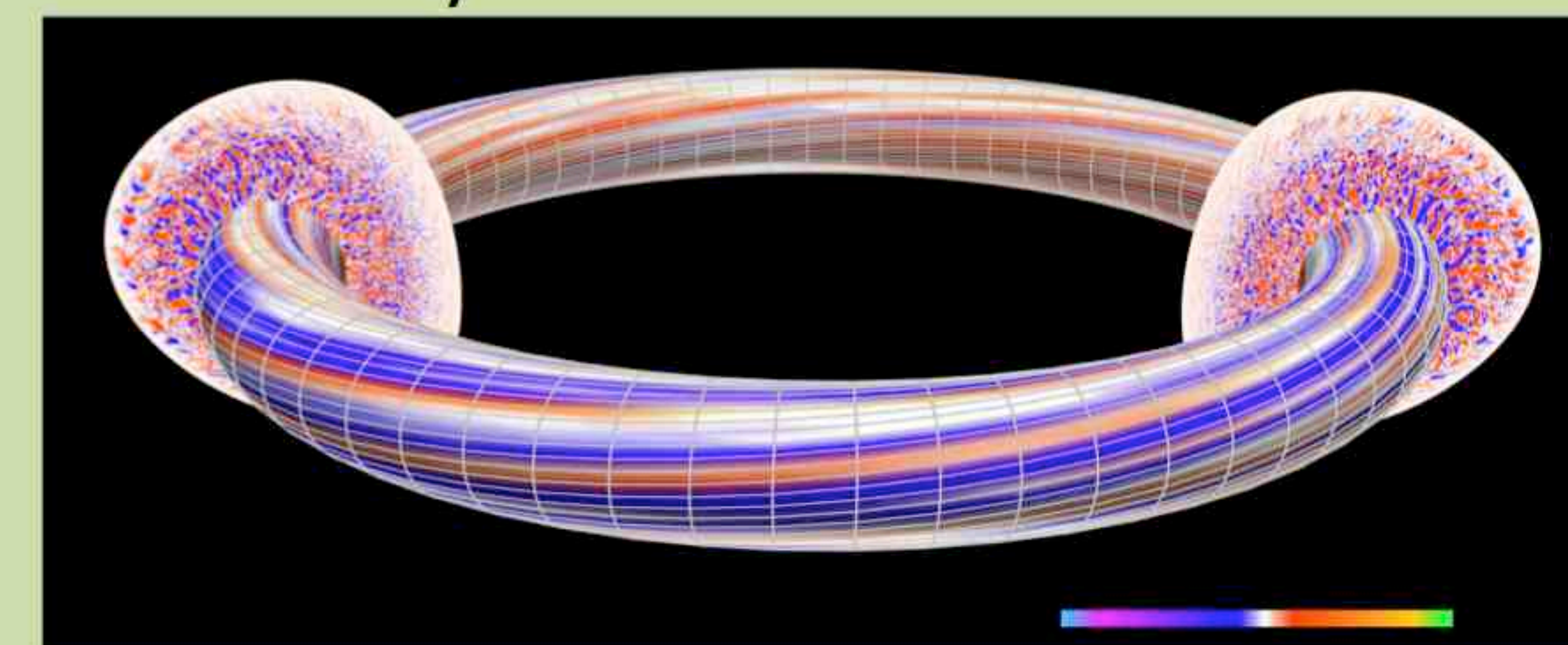
```

## SHIFT ALGORITHMS: TWO-SIDED (MPI-1) VS. ONE-SIDED (PGAS)

- MPI algorithms have been extensively researched and optimized in the past
- Multithreaded MPI algorithms exploiting shared memory work-sharing constructs + novel overlapping techniques, e.g., using OpenMP task, have been developed
- All optimized MPI algorithms share the large bulk data transfers to exchange moving particles following rule that performance is optimized by sending fewer & larger messages
- Vendor-supplied MPI library that has been highly tuned for the XE6/Gemini interconnect
- Exploit the one-sided nature of the PGAS model
- Lightweight data transfers without synchronization between the sending & receiving images
- Sending more frequent smaller messages enables the Coarray approach to outperform the MPI-1 implementations due to enhanced communication and computation overlap as well as the better network bandwidth utilization
- A similar strategy in MPI-1 cannot be implemented and very hard in MPI-2
- One-sided messaging semantics as language constructs (Coarrays in Fortran 2008) allow compilers to directly reference remote memory and to apply communication optimizations
- HW support for PGAS, as provided by the recent Cray XE6 Gemini interconnect, is essential to realize the performance potential; but, still good performance improvements on XT4!



## APPLICATION: Gyrokinetic Tokamak Simulation (GTS)

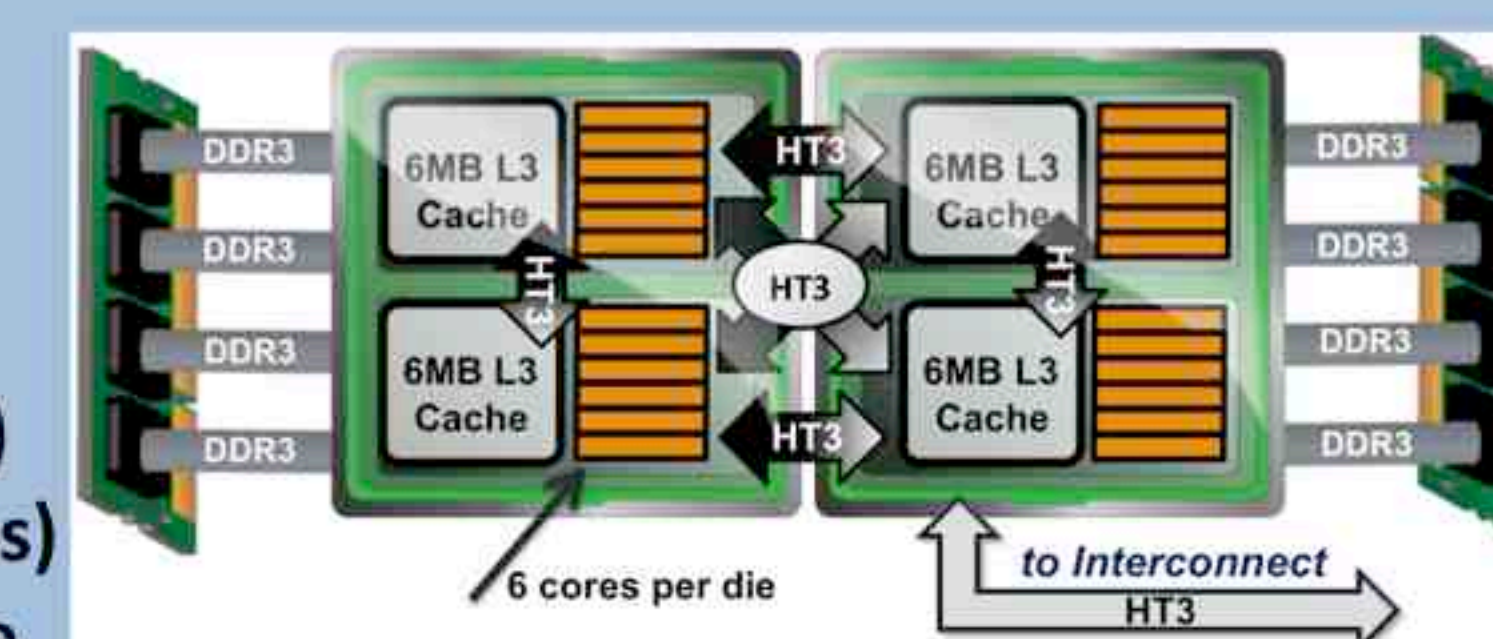


- A global 3D Particle-In-Cell (PIC) code with MPI & OpenMP support
- Developed to study plasma microturbulence in toroidal, magnetic confinement devices called tokamaks
- Microturbulence is a complex, nonlinear phenomenon that is believed to play a key role in the confinement of energy and particles in fusion plasmas
- Due to one of the levels implemented for parallelism in GTS, particles from one toroidal domain to another while they travel around the torus
- Shift phase** represents the most communication intense routine in GTS and will gain in importance when scaling GTS to Petascale or even Exascale supercomputers.

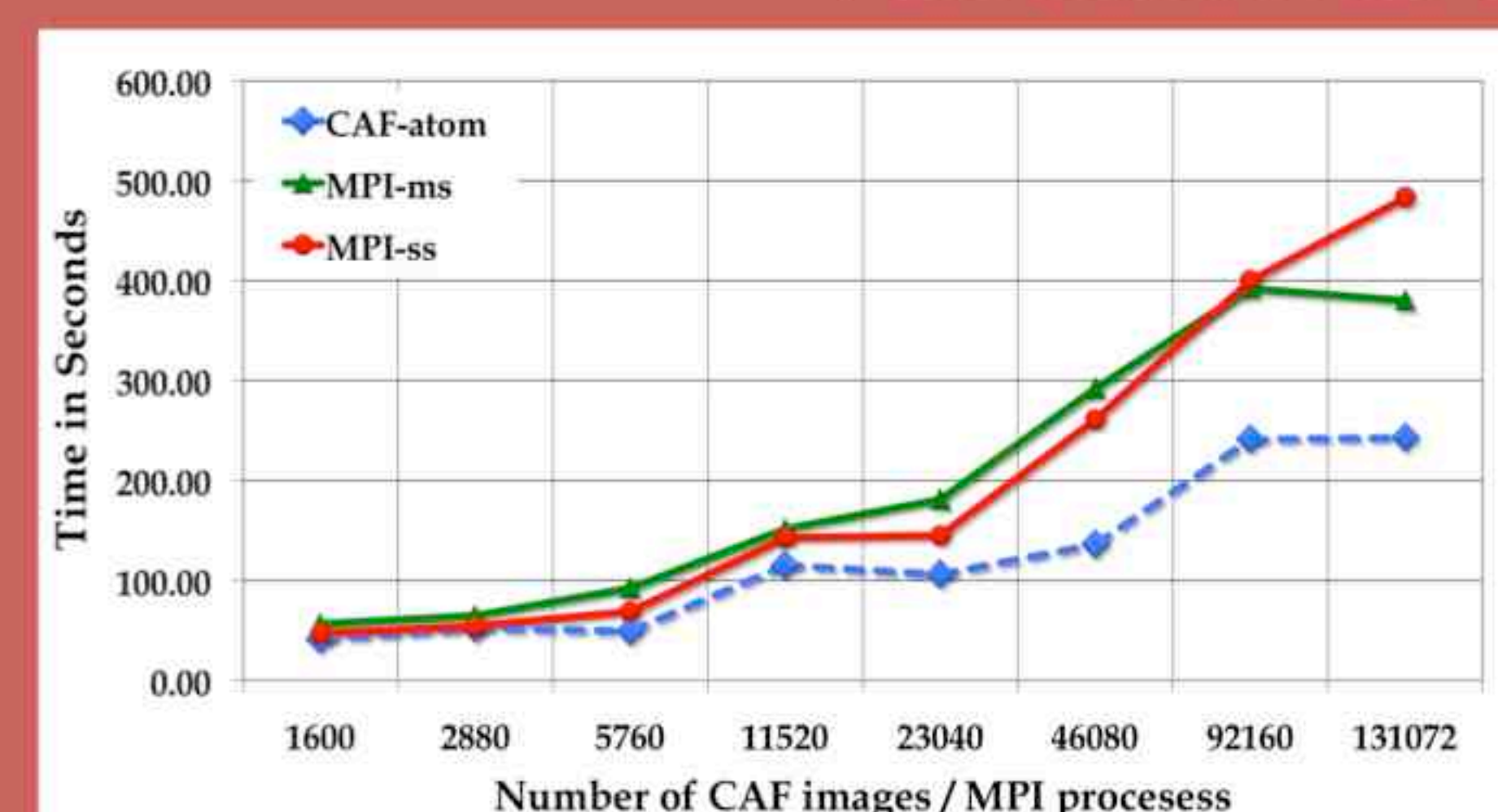
## ARCHITECTURE: NERSC CRAY XE6 "HOPPER"

### Compute Node Configuration:

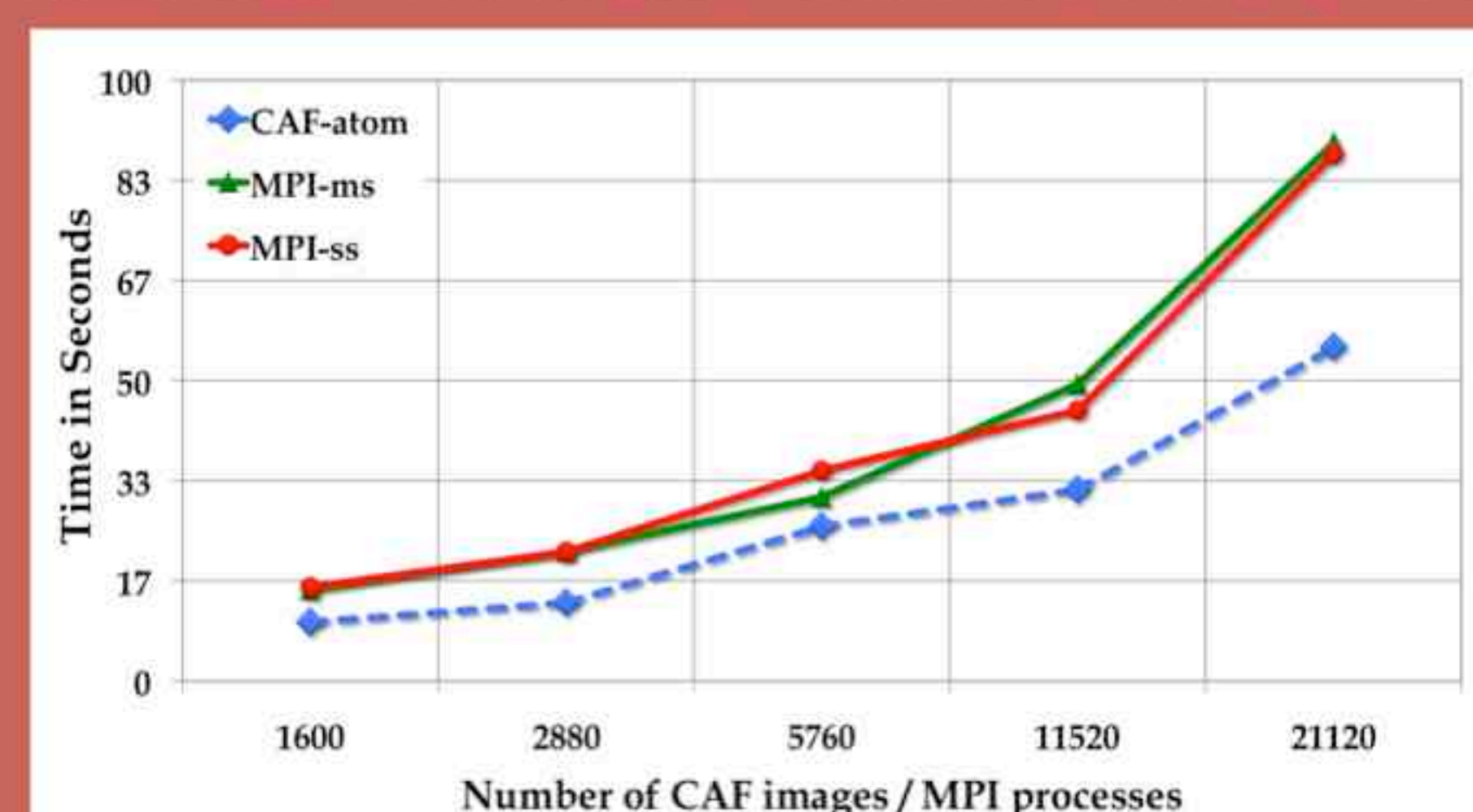
- 6384 nodes
- 2 twelve-core AMD 'MagnaCours' 2.1 GHz processors per node (NUMA)
- 24 cores per node (153,216 total cores)
- 1.28 Peta-flops for the entire machine
- Each core has their own L1 and L2 caches, with 64 KB and 512KB respectively
- 6 MB L3 cache shared between 6 cores on the Magna-Cours processor
- 4 DDR3 1333 MHz memory channels per twelve-core 'MagnaCours' processor



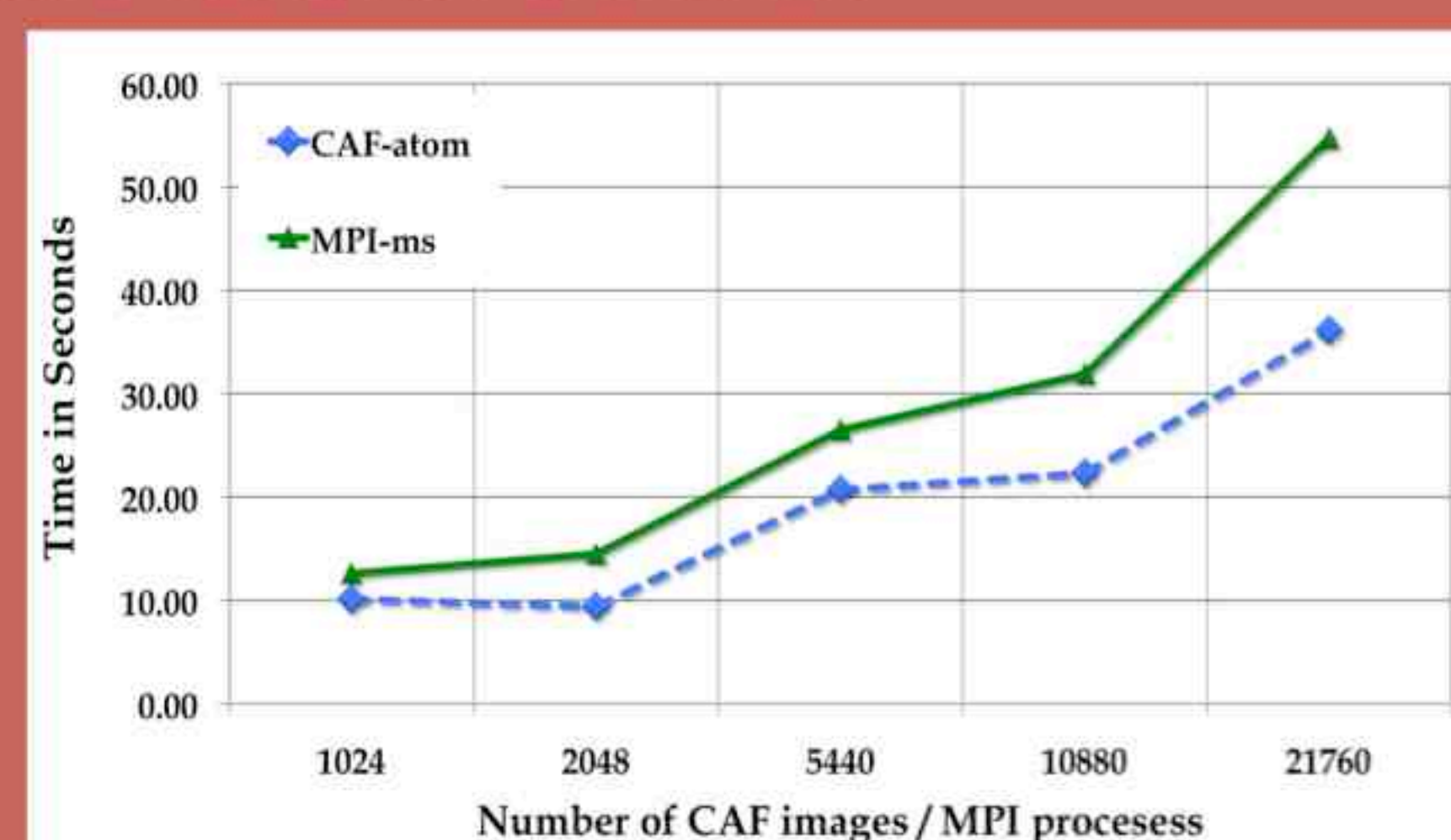
## PERFORMANCE: EVALUATED IN A BENCHMARK SUITE & IN THE REAL APPLICATION



Single-Threaded (Benchmark Suite)



Multi-Threaded (Benchmark Suite)



Multi-Threaded (GTS)

### Gemini Characteristics:

- 168 GB/sec routing capacity
- Scales to over 100,000 network endpoints
- Link-level reliability and adaptive routing
- Provides global address space
- Efficiently supports MPI, one-sided MPI, Shmem, UPC, Co-Array Fortran
- Supports millions of MPI messages per second
- Internode latency on the order of 1 microsecond
- Bandwidth of 9.8 GB/sec per Gemini chip

